

Noise-contrastive estimation of normalising constants and GANs

Fonctions génériques

```
library(ggplot2)
library(reshape)
library(matrixStats)
```

Algorithme d'Hasting

Utilité : simuler selon $p_m(., \psi)$ pour un paramètre ψ choisi.

Argument	Type	Exemple	Indication
x	vecteur	rcauchy(100, 0, 1)	notre échantillon de densité inconnue
n	entier	100	taille de la simulation
psi	vecteur	c(0,1)	paramètres de la fonction h
h	fonction		fonction qui retourne $\overline{p_m}(., \psi)$

```
hasting = function(x, n, psi, h){
  y = c()
  y = append(y, sample(sample(x, 1)))
  for (i in 2:n){
    y_ = rnorm(1, y[i-1], 1)
    u = runif(1)
    if ( u <=
          (h(y_,psi) * dnorm(y_, y[i-1], 1))
          / (h(y[i-1],psi) * dnorm(y[i-1], y_, 1))
    ){
      y = append(y, y_)
    }
    else {
      y = append(y, y[i-1])
    }
  }
  return (y)
}
```

Note : on peut très certainement écrire sous forme matricielle cette fonction pour une meilleure performance.

MC MLE

Utilité : retourne une estimation des paramètres selon la méthode décrite dans le papier de Geyer.

```
mc_mle = function(x, psi, h){

  m = length(x)

  y = hasting(x, m, psi, h)
```

```

L = function(theta){
  return(sum(log(h(x,theta)/h(x,psi))) - m*log(mean(h(y,theta)/h(y,psi))))
}

theta = optim(
  par = rep(1,length(psi)),
  gr = "CG",
  control = list(fnscale=-1),
  fn = L
)$par

return(theta)
}

```

NCE

Utilité : Retourne l'estimation de la constante et des paramètres.

Argument	Type	Exemple	Indication
x	vecteur	rcauchy(100, 0, 1)	notre échantillon de densité inconnue
law_y	fonction	rnorm	fonction qui retourne un échantillon suivant la loi p_n
n	entier	100	taille de l'échantillon de bruit suivant la loi p_n
params_y	vecteur	c(0,1)	arguments de la fonction law_y
log_pm	fonction		fonction qui retourne le logarithme de la densité p_m
log_pn	fonction		fonction qui retourne le logarithme de la densité p_n
size_theta	entier	3	taille de θ , vaut habituellement 2 ou 3
method	string	"CG"	méthode d'optimisation, habituellement "CG" ou "BFGS"

```

nce = function(x, law_y, params_y, log_pm, log_pn, size_theta, n, methode = "CG"){

  y = do.call(law_y, c(list(n),params_y))

  m = length(x)

  h = function(u, theta){
    return( 1 / (1 + n/m * exp(log_pn(u) - log_pm(u, theta))))
  }

  J = function(theta){
    return( sum(log(h(x, theta))) + sum(log(1 - h(y, theta))) )
  }

  theta = optim(
    par = rep(1, size_theta),
    gr = methode,
    control = list(fnscale=-1),
    fn = J
  )$par

  return(c(theta[-size_theta], exp(-theta[size_theta])))
}

```

Graphiques

Utilité : afficher l'histogramme pour un échantillon de données x .

```
print_hist = function(x) {  
  df = data.frame(x = x)  
  hist_x = ggplot(df, aes(x=x)) + geom_histogram(aes(y = stat(count) / sum(count)), bins = 20, color="white", fill="black")  
  print(hist_x)  
}
```

Utilité : pour NCE, afficher l'évolution des paramètres au fur et à mesure de l'augmentation de n (la dimension de l'échantillon de bruit)

```
NCE_evol_params = function(x, law_y, params_y, log_pm, log_pn, size_theta, ratio, steps, labels, methode) {  
  
  # Creation de l'abscisse  
  m = length(x)  
  N = seq(0, m*ratio, length.out = steps + 1)  
  
  # Creation de l'ordonnée  
  theta = c()  
  for (n in N) {  
    theta = append(theta, nce(x, law_y, params_y, log_pm, log_pn, size_theta, n, methode))  
  }  
  
  # Formatage des données  
  theta = t(rbind(matrix(theta, nrow = size_theta), N))  
  df = as.data.frame(theta)  
  df_melted = melt(df, id.vars = "N")  
  
  # Plot  
  plot_df = ggplot(df_melted, aes(x = N, y = value)) +  
    geom_line(aes(color = variable, group = variable)) +  
    geom_point(aes(color = variable, group = variable)) +  
    labs(title = "Evolution des paramètres par rapport au bruit", x = "n (taille du bruit)", y = "Paramètres")  
    scale_color_manual(labels = labels, values = c("blue", "red", "orange"))  
  
  print(plot_df)  
  
  return(theta)  
}
```

Utilité : utilise le bootstrap sur x pour estimer les paramètres

```
NCE_bootstrap = function(x, law_y, params_y, log_pm, log_pn, size_theta, n, nb_of_boot, labels, methode) {  
  
  theta_bootstrap = c()  
  for (i in 1:nb_of_boot) {  
    x_bootstrap = sample(x, size = m, replace=TRUE)  
    theta_bootstrap = append(theta_bootstrap, nce(x_bootstrap, law_y, params_y, log_pm, log_pn, size_theta, n, methode))  
  }  
  
  # Formatage des données  
  theta_bootstrap = matrix(theta_bootstrap, nrow = size_theta)  
  
  # Plot
```

```

size_bootstrap = 1:nb_of_boot
df = as.data.frame(cbind(t(rowCumsums(theta_bootstrap)) / size_bootstrap, size_bootstrap))
df_melted = melt(df, id.vars = "size_bootstrap")

plot_df = ggplot(df_melted, aes(x = size_bootstrap, y = value)) +
  geom_line(aes(color = variable, group = variable)) +
  labs(title = "Evolution des paramètres par bootstrap", x = "Taille du bootstrap", y = "Paramètres", color = "variable")
  scale_color_manual(labels = labels, values = c("blue", "red", "orange"))

print(plot_df)

return(rowMeans(theta_bootstrap))
}

```

Note : il faudrait optimiser le temps de calcul de ces fonctions, peut-être en matriciel au lieu des boucles ou bien avec du calcul en parallèle sur CPU/GPU

Exemple basique : la loi normale

Soit x l'échantillon de taille m obtenu selon la loi de densité inconnue p_d .

On considère ici que p_d appartient à la famille de fonctions paramétrées par $\theta = (c, \mu, \sigma)$ suivante :

$$p_m(u; \theta) = \frac{1}{Z(\mu, \sigma)} \times \exp\left[-\frac{1}{2}\left(\frac{u - \mu}{\sigma}\right)^2\right] \quad \text{d'où} \quad \ln(p_m(u; \theta)) = c - \frac{1}{2}\left(\frac{u}{\sigma} - \frac{\mu}{\sigma}\right)^2$$

```

pm_barre = function(u, theta){
  return(exp(-0.5 * ((u - theta[1]) / theta[2]) ** 2))
}

log_pm = function(u, theta){
  return(theta[3] - 1/2 * (u/theta[2] - theta[1]/theta[2]) ** 2)
  # theta[1] = mu / theta[2] = sigma / theta[3] = c
}

log_pn_cauchy = function(u){
  return(log(dcauchy(u, mean(x), sd(x))))
}

m = 10000
n = 100000
x = rnorm(m, 2, 4)
size_theta = 3

# METHODE MC MLE
mc_mle(x, c(mean(x), sd(x)), pm_barre)

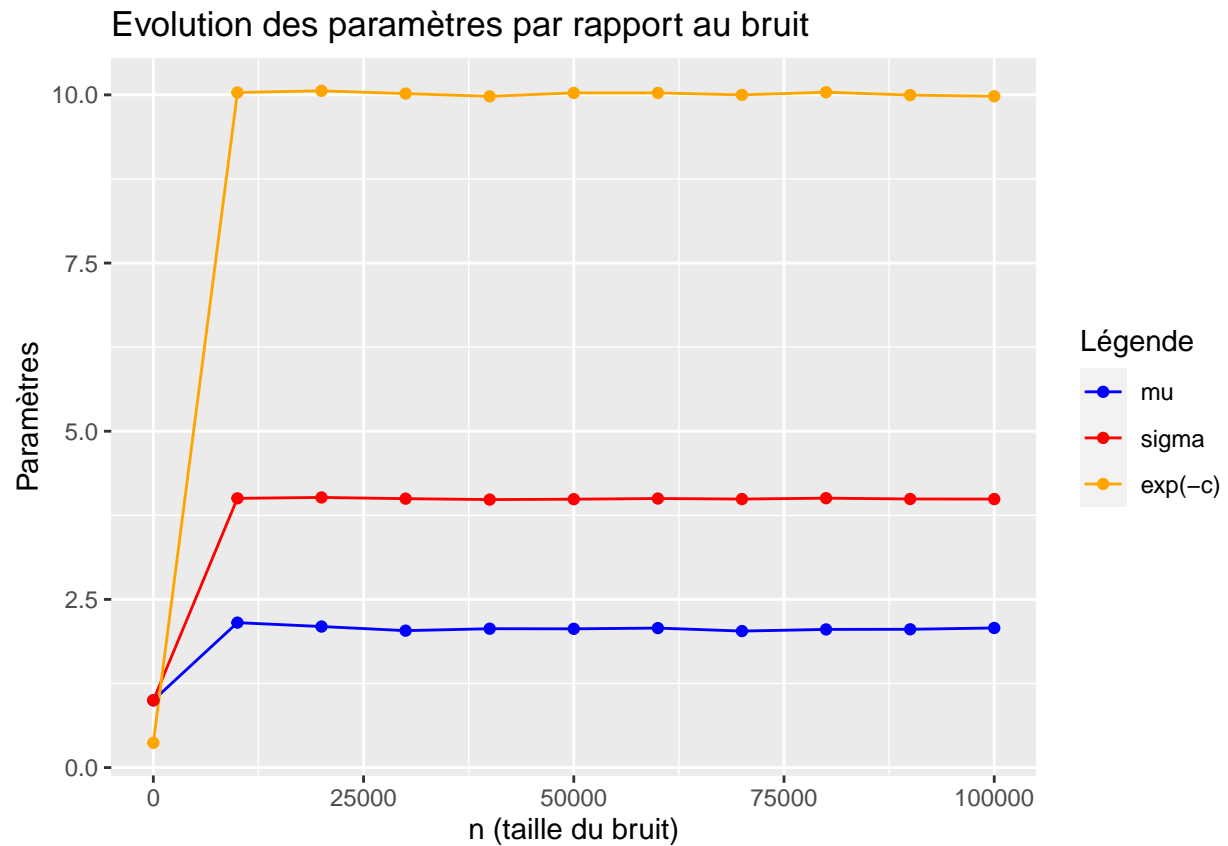
## [1] 1.408763 3.868717

# METHODE GEYER
nce(x, rcauchy, c(mean(x), sd(x)), log_pm, log_pn_cauchy, size_theta, n)

## [1] 2.057044 3.984856 9.998000

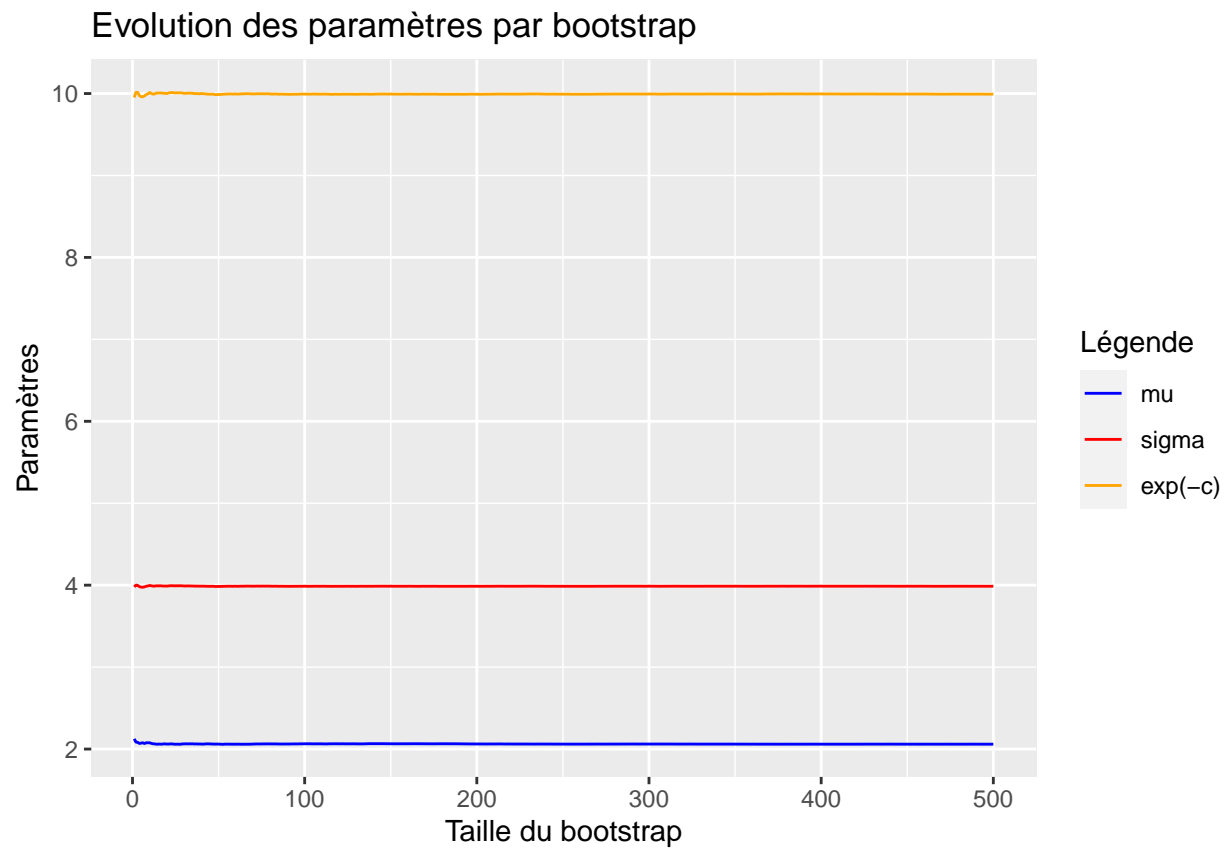
```

```
NCEevol_params(x, rcauchy, c(mean(x),sd(x)), log_pm, log_pn_cauchy, size_theta, 10, 10, c("mu", "sigma", "exp(-c)"))
```



```
##
##      N
## [1,] 1.000000 1.000000 0.3678794 0e+00
## [2,] 2.154383 4.002365 10.0351628 1e+04
## [3,] 2.096231 4.015271 10.0596207 2e+04
## [4,] 2.035205 3.997363 10.0187334 3e+04
## [5,] 2.064043 3.983743 9.9781736 4e+04
## [6,] 2.062410 3.989303 10.0292124 5e+04
## [7,] 2.073276 3.999511 10.0289329 6e+04
## [8,] 2.029126 3.991589 9.9998321 7e+04
## [9,] 2.053478 4.004637 10.0390035 8e+04
## [10,] 2.055587 3.992300 9.9971635 9e+04
## [11,] 2.075100 3.991190 9.9792127 1e+05
```

```
(NCE_bootstrap(x, rcauchy, c(mean(x),sd(x)), log_pm, log_pn_cauchy, size_theta, n, 500, c("mu", "sigma", "exp(-c)"))
```



```
## [1] 2.059720 3.987117 9.992494
```